

The importance of directed triangles with reciprocity: patterns and algorithms*

C. Seshadhri
Sandia National Laboratories[†]
Livermore, CA
scomand@sandia.gov

Ali Pinar
Sandia National Laboratories
Livermore, CA
apinar@sandia.gov

Nurcan Durak
Sandia National Laboratories
Livermore, CA
nurcan.durak@gmail.com

Tamara G. Kolda
Sandia National Laboratories
Livermore, CA
tgkolda@sandia.gov

ABSTRACT

The computation and study of triangles in graphs is a standard tool in the analysis of real-world networks. Yet most of this work focuses on undirected graphs. Real-world networks are often directed and have a significant fraction of reciprocal edges. While there is much focus on directed triadic patterns in the social sciences community, most data mining and graph analysis studies ignore direction.

But how to we make sense of this complex directed structure? We propose a collection of *directed closure values* that are analogues of the classic *transitivity* measure (the fraction of wedges that participate in triangles). We perform an extensive set of triadic measurements on a variety of massive real-world networks. Our study of these values reveal a wealth of information of the nature of direction. For instance, we immediately see the importance of reciprocal edges in forming triangles and can measure the power of transitivity. Surprisingly, the chance that a wedge is closed depends heavily on its directed structure. We also observe striking similarities between the triadic closure patterns of different web and social networks.

Together with these observations, we also present the first sampling based algorithm for fast estimation of directed triangles. Previous estimation methods were targeted towards undirected triangles and could not be extended to directed graphs. Our method, based on wedge sampling, gives orders of magnitude speedup over state of the art enumeration.

[†]Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.

*This work was funded by the GRAPHS Program at DARPA and the Laboratory Directed Research and Development program at Sandia National Laboratories.

Categories and Subject Descriptors

G.2.2 [Discrete Mathematics]: Applications

General Terms

Algorithms, Theory

Keywords

wedge closure, transitivity, directed triangles, sampling

1. INTRODUCTION

The study of triangles is by now a classic tool in the analysis of large-scale networks. The focus on triangles has its roots in a variety of disciplines: in social sciences as a manifestation of various theories, in physics as local measures of clustering, in biology as *motifs*. Yet most contemporary data mining and massive graph analysis first convert real-world interaction data (think of this as a graph with attributes) into an undirected graph, and *then* work on this graph. This is a very fruitful method, since the complexity of the underlying problem is reduced, and we still get a significant amount of information. Nonetheless, it is a major challenge to account for the attributes on edges.

The most common attribute for edges is *direction*. For example, most social networks, web networks, and product networks are all truly directed networks. Moreover, directed networks often have a significant percentage of *reciprocal edges*. Newman et al. [18] shows that the fraction of such edges in commonly studied graphs is quite high, and subsequent studies underlined the importance of such edges in virus/news spreading and understanding the network formation [10, 17, 14].

The set of triangles (and wedges) involving directed and reciprocal edges is rich and holds information about the underlying dynamics [13, 16, 9, 8, 23]. But it is challenging to make sense of this information and also compare different graphs (from varied sources) along these metrics. Furthermore, computation of triangles becomes quite expensive for large graphs.

1.1 Some preliminaries

We focus on a directed graph (digraph) $G = (V, E)$. In a standard digraph, all edges are just ordered pairs of vertices of the form (i, j) . We will think of the graph as having two

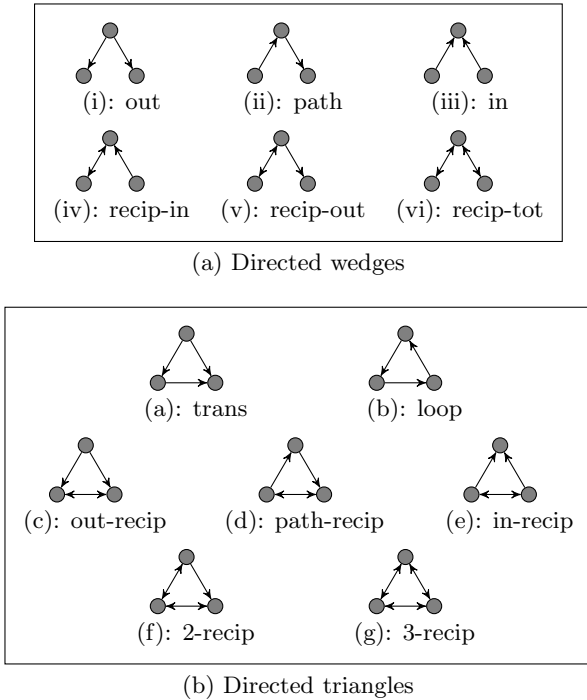


Figure 1: Directed structures

different types of edges: basic and reciprocal. A reciprocal edge is technically a pair $\{(i, j), (j, i)\}$, which we merge into a single reciprocal edge. *We do not think of a reciprocal edge as containing two directed edges, but consider it to be a special edge on its own.* In our figures, reciprocal edges are depicted as double-headed arrows. We define reciprocity of a graph, r , as the ratio of the number of reciprocal edges to the total number of edges. Note that our definitions lightly differ than that of [18].

A *wedge* is a pair of edges that share an endpoint, and a *triangle* is a set of three (unparallel) edges that are incident on a set of three vertices. We have 6 different types of wedges and 7 different types of triangles. We give more details about these structures in §2. We give the list of directed wedges and triangles with reciprocity in Figure 1. The earliest construction of this list is by Holland and Leinhardt [13].

1.2 Main results of this paper

- **Definition of directed closures:** We generalize the classic notion of *transitivity* (pg. 243 of [25]), which is also called the global clustering coefficient, to digraphs. This leads to a set of 15 closure values that provide a triadic pattern of a digraph.

- **Computation of directed triangles and closures:** We extend the basic method of *wedge sampling* [19, 20] to approximating the counts of directed triangles. This gives fast scalable algorithms with provable error bounds. While triangle counting is a well-studied problem, we present the first algorithm that works for digraphs. This algorithm enables efficient computation of all closure values.

We perform experiments on a set of publicly available datasets. We present the directed closure information in a succinct form that allows comparison of different graphs. This leads to a series of observations.

- **Heterogeneity of closure:** We find the closure fractions of wedges vary greatly depending on the wedge type. In-wedges ((iii) in Figure 1a) usually dominate the graph but are rarely closed. In many cases, *all other* wedge types close frequently.

- **Reciprocity induces closure:** For almost every graph we analyze, the presence of a reciprocal edge in a wedge greatly increases the chance of closure. In other words, wedges with reciprocal edges participate in triangles more frequently than (uniform) random wedges.

- **The power of transitivity:** Loops and path-recip triangles ((b) and (d) in Figure 1b) are very infrequent. These triangles contain a transitive wedge that is not reciprocated, and the fact that they are so rare suggest the power of transitivity in the underlying dynamics. This appears to validate the importance of transitivity, as posited by Holland and Leinhardt [13] in the social science community (Recent results of Leskovec et al. [15] in signed networks make comparable observations). These observations also underscore the importance of reciprocity, since this distinguishes triangles without transitivity from those that have it.

- **The non-randomness of direction:** We define a simple random model of direction in an underlying undirected graph and compute directed closure values for this model. The predictions from this are significantly different from the actual data, showing that our findings indicate a deep directed structure in real-world networks.

What is the significance of these results? First, we feel that these observations show the importance of direction and reciprocity, which we believe is not emphasized enough in analyses of social networks. Designing meaningful measures related to directed triangles and interpretable presentations is an important step in understanding digraphs. We also need efficient algorithms to compute such measures. We hope that this work is a step in this process. The wealth of information that is obtained by looking at directed closure values (at least in the authors' opinion) shows the importance of the directed closure values.

These values also inform *graph modeling* because they provide formal measures that models can be tested against. It has been observed before that existing graph models have little to no reciprocity [7], so no model can even come close to matching directed closures. We have no models that even come close to recreating structure of digraphs. This is probably a very difficult problem, but greater insight into directed closures might help in making progress.

An attribute of networks that we ignore is *sign* (a positive versus negative relationship). Many social science theories focus on sign in networks, and recent work by Leskovec et al. [15] studies signed networks. It would be interesting to extend our work to signed *and* directed networks.

1.3 Previous work

The earliest study of directed triads with reciprocity, to our knowledge, is in the social sciences, by Holland and Leinhardt [13]. They explicitly list the 16 different possible triads (including the 3 patterns with at most one edge) and count them in various social networks of the time. They also try to measure the effects of reciprocity in network formation. This is called the *triad census*. Skvoretz [21] and Skvoretz et al. [22] use these numbers of predict various biases in network formation. In a more recent study, Faust [9] computes a triad census on many graphs to compare their structure.

Most of this work has been restricted to small data sets (at most hundreds of nodes). Finding such triads has been referred to as *motif finding* in the bioinformatics community [16]. Simpler versions of triad census counts have also been used to analyze gaming data [23].

A classic *local* measure of triangle density is *clustering coefficient*, introduced by Watts and Strogatz [26]. Fagiolo [8] proposes a local clustering coefficient measure for directed networks, though he ignores reciprocity. Ahnert and Fink [2] construct “clustering coefficients signatures” from these measures and classify directed networks.

Leskovec et al [15] study *signed* networks and validate (and extend) the *theory of balance* [11, 3]. They study the behavior of signed triangles to show that theory of balance does not suffice to explain networks. They also look at direction, but their datasets do not involve much reciprocity. It would be interesting to combine their work with our measures of directed closures.

2. THE DIRECTED CLOSURES

We begin with some notation and introduction to the directed structures in Figure 1a and Figure 1b. We use small Roman numerals to index the types of wedges, and small Latin letters for triangles. Furthermore, ψ is used to denote a variable wedge type, and τ for a variable triangle type. We also give some names for further reference. (Holland and Leinhardt [13] have a naming scheme for directed triads that involve a triple of numbers with a letter. We deviate from this notation because its easier to remember names than 3 digit codes.)

We stress that these types form a partition of all wedges and triangles. Since reciprocal edges are special, we do not think of (say) the recip-out wedge containing an out wedge.

For each vertex v , we have three associated degrees: the indegree, outdegree, and reciprocal degree. These are denoted by d_v^{\leftarrow} , d_v^{\rightarrow} , and d_v^{\leftrightarrow} . The total degree $d_v = d_v^{\leftarrow} + d_v^{\rightarrow} + d_v^{\leftrightarrow}$. We mention some of the salient features of these directed structures.

- **Basic vs reciprocal structures:** The structures without reciprocal edges form the first rows in both Figure 1a and Figure 1b. There are only 3 types of wedges and 2 types of triangles, underscoring the importance of reciprocity.

- **Cyclic relations:** Triangle types (b), (d), (f), (g) all contain a cycle, and there is a progression of 0, 1, 2, and 3 reciprocal edges.

- **The table of $\chi(\psi, \tau)$ values:** Different triangle types naturally contain different types of wedges. This information is summarized by the function $\chi(\psi, \tau)$, which we define as the number of type ψ wedges in type τ triangles. The list of nonzero values of $\chi(\psi, \tau)$ is provided in Table 1. Each row contains the wedge information of that triangle type. There are 15 nonzero entries in this table.

- **Wedge counts:** For vertex v , let $W_{v,\psi}$ be the set of ψ -wedges centered at v . It is routine to compute $|W_{v,\psi}|$ given the degrees of v . This is summarized in Table 2.

2.1 (ψ, τ) -closure

The transitivity (or global clustering coefficient) is defined as $3|T|/|W|$ (T is the set of triangles and W is the set of wedges). Semantically, this is the fraction of wedges that participate in triangles.

In the undirected setting, a wedge is called *closed* if it participates in a triangle and open otherwise. We say that

		Wedge types (ψ)					
Triangle types (τ)		i	ii	iii	iv	v	vi
	a	1	1	1			
	b		3				
	c	1				2	
	d		1		1	1	
	e			1	2		
	f				1	1	1
	g						3

Table 1: Number of occurrences of each wedge type per triangle type: $\chi(\psi, \tau)$.

ψ						
$W_{v,\psi}$	$\binom{d_v^{\rightarrow}}{2}$	$d_v^{\leftarrow} d_v^{\rightarrow}$	$\binom{d_v^{\leftrightarrow}}{2}$	$d_v^{\leftarrow} d_v^{\leftrightarrow}$	$d_v^{\rightarrow} d_v^{\leftrightarrow}$	$\binom{d_v^{\leftrightarrow}}{2}$

Table 2: Number of wedges per vertex for each wedge type.

a ψ -wedge is τ -closed if the wedge participates in a type τ triangle. The (ψ, τ) -closure, $\kappa_{\psi,\tau}$, is the fraction of ψ -wedges that are τ -closed. Formally, let W_{ψ} be the set of ψ -wedges and T_{τ} be the set of τ -triangles.

$$\kappa_{\psi,\tau} = \frac{\chi(\psi, \tau)|T_{\tau}|}{|W_{\psi}|}$$

The number of ψ -wedge in τ -triangles is $\chi(\psi, \tau)|T_{\tau}|$. Note that if type τ triangles contain no type ψ wedge, then this quantity is just zero because of $\chi(\psi, \tau)$. As mentioned earlier, there are 15 non-trivial (ψ, τ) -closures.

2.2 Representations

We create a *directed closure chart*, which combines all the $\kappa_{\psi,\tau}$ values. We give an example for the web-Google [27] graph in Figure 2. The bars on the x -axis are indexed by the different wedge types, and the y -axis is $\kappa_{\psi,\tau}$. We make a stacked bar chart with the different closure values, where the triangle types are shown in 7 different colors. For example, the blue part of the first bar is the fraction of out-wedge closing into trans triangles ($\kappa_{i,a}$). Some of the salient features:

1. **Single closure value:** Consider some wedge type and triangle type (say out-wedge and trans-triangle). The value $\kappa_{i,a}$ is shown by the height of the blue part of the first bar. The height of the blue part in the second bar show the fraction of path-wedges that are closed into a trans-triangle.

2. **Total closure of wedge type:** The total height of the bar is total fraction of closed wedges of that type. For example, we see that in-wedges close infrequently.

3. **Percentage of wedge type:** Underneath the wedge pictures is the percentage of that wedge type.

4. **Percentage of triangle type:** Underneath the legend for triangles is the percentage of that triangle type.

5. **Undirected transitivity:** The value of κ is marked by a thick dashed line.

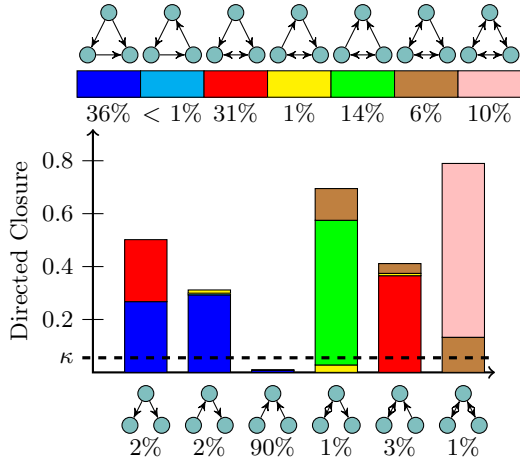


Figure 2: Directed closure for web-Google

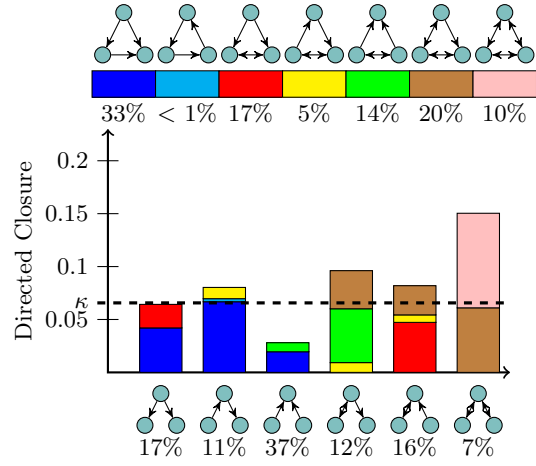


Figure 5: Directed closure for soc-Epinions1

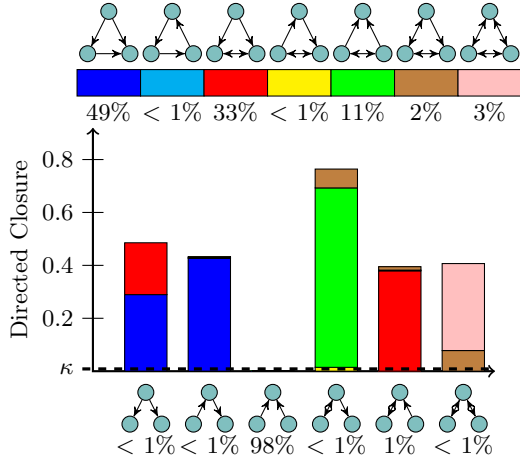


Figure 3: Directed closure for web-Stanford

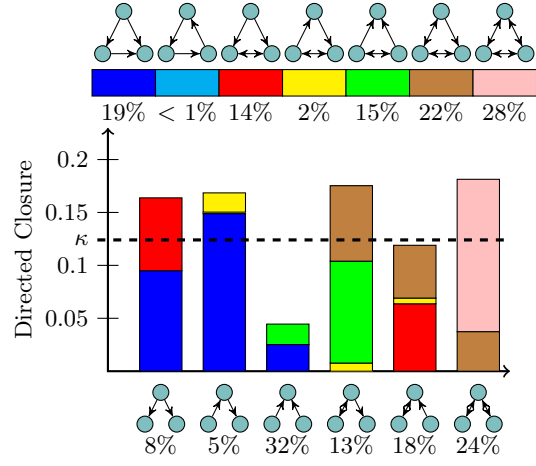


Figure 6: Directed closure for livejournal

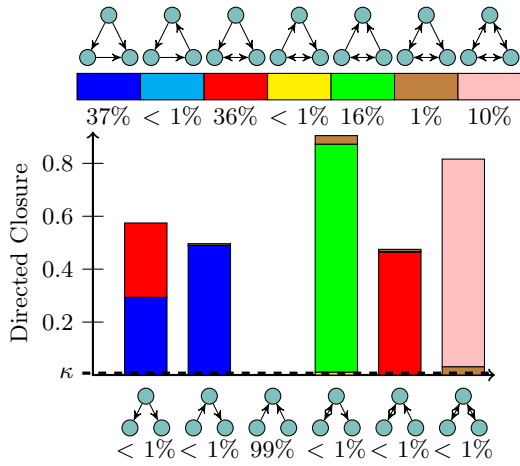


Figure 4: Directed closure for web-BerkStan

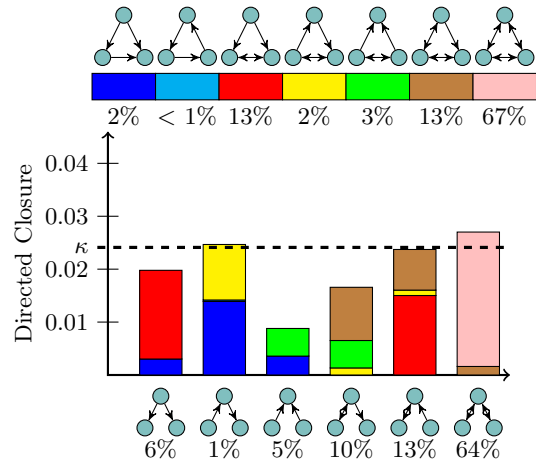


Figure 7: Directed closure for soc-Slashdot0902

Table 3: Properties of the graphs

Graph Name	V	E	W	T	r	κ
amazon0505	410K	3357K	73M	3951K	0.55	0.162
soc-Slashdot0902	82K	870K	75M	603K	0.84	0.024
web-Stanford	282K	2312K	3944M	11330K	0.28	0.009
web-BerkStan	685K	7601K	27983M	64691K	0.25	0.007
wiki-Talk	2394K	5021K	12594M	9204K	0.14	0.002
web-Google	876K	5105K	727M	13392K	0.31	0.055
soc-Epinions1	76K	509K	74M	1624K	0.41	0.066
web-NotreDame	326K	1470K	305M	8910K	0.52	0.088
youtube-links	1158K	4945K	1474M	3057K	0.79	0.006
flickr-links	1861K	22614K	14675M	548659K	0.62	0.112
soc-livejournal	5284K	76938K	7519M	310877K	0.73	0.124

3. OBSERVATIONS ON CLOSURE CHARTS

We analyze the directed closure properties of various real graphs, whose properties are presented in Table 3. In this table, $|V|$, $|E|$, $|W|$, and $|T|$ correspond to the number of vertices, edges, wedges, and triangles, respectively. The *reciprocity* r is the fraction of total edges that are reciprocal edges. The undirected transitivity ($3|T|/|W|$) is given by κ .

3.1 Similarities of directed closures

Figure 2, Figure 3, and Figure 4 have the closure charts for three different web graphs: web-Google, web-Stanford, and web-BerkStan [27]. These graphs have vertices for web-pages and directed edges for web links. Figure 5, Figure 6, and Figure 7 have the charts for three social networks [27]. The vertices of soc-Epinions are member of Epinions, a consumer review site. A directed edge between users shows a trust relationship originating from one user (these are signed by trust/distrust, which we ignore). The vertices of soc-Slashdot [27] are users and edges represent tagging as friend or foe. The vertices of soc-livejournal [5, 1] are Slashdot users with edges denoting friendship (which is one-way).

Observe the uncanny similarity of the closure charts web graphs, despite them being from different sources (and different sizes). The color patterns are remarkably similar, showing similar distributions of different closures. The social networks show more variation, but the overall structure of the charts is not far from the web graphs. In general, we note that in-wedges rarely close and reciprocal wedges close much more frequently.

3.2 Heterogeneity of closure

The heterogeneity of wedge closure is quite clear from all the closure charts. Focus on the web graphs. Other than in-wedges, all other wedge types close (quite) frequently. The undirected transitivity is always below 0.05, but specific wedge types close more than 50% of the time (shown by the total height of the bar). In-wedges form a dominant majority of all wedges (more than 98%) but close infrequently. Indeed, the low value of transitivity is explained by the high percentage yet low closure of in-wedges.

The picture is not as dramatic in the social networks, but there is some variation in closures over the wedge types. Quite consistently, in-wedges do not close and recip-tot-wedges close more frequently.

3.3 Effect of reciprocity on closure rates

How does reciprocity change the chance of closure? Observe that in, path, and out-wedges contain no reciprocal

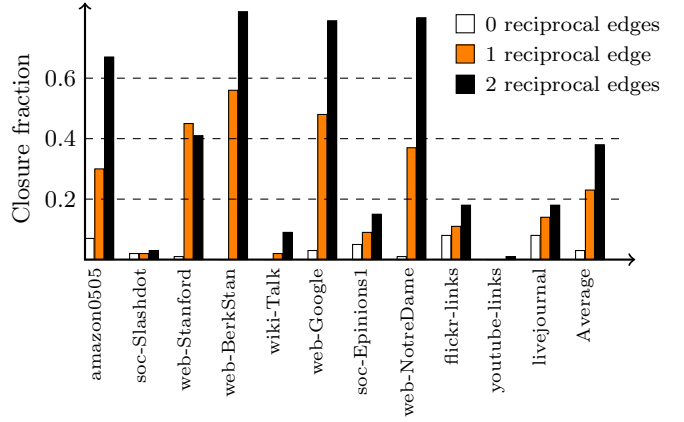


Figure 8: How reciprocity increases closure rates: the x-axis gives over various graphs. The colored bars correspond to wedges with 0, 1, or 2 reciprocal edges. The y-axis gives the fraction of those wedges that close (into any triangle).

edges, recip-in and recip-out-wedges have exactly 1 reciprocal edge, and recip-tot has 2 reciprocal edges. As the charts clearly indicate, having reciprocal edges increases the chance of closure a wedge. We do a comprehensive calculation on a variety of graphs in Figure 8.

Consider a graph and choose k from $\{0, 1, 2\}$. Fix the set of wedges with k reciprocal edges, and look at the fraction of those that close (into any triangle). This gives the data presented in Figure 8. Observe how there is consistently a monotonic (and often dramatic) increase in closure fractions as reciprocity increases. The average of chance of closure for a wedge without reciprocal edges is only 3%. But this number goes to 23% if one of the edges is reciprocal and further increases to 38% when both edges are reciprocal. This finding is consistent with the earlier reports about reciprocal edges, indicating stronger ties between two vertices [18, 10, 17, 14]. It also underscores how important it is to consider direction in networks.

3.4 The power of transitivity

Throughout the closure charts, one notices in infrequency of loop and path-recip-triangles. These are colored light blue and yellow, and one can see how little of those colors are present (or one can directly look at their percentages). Let us focus on triangles that contain a cycle showing a “cyclic” relationship. These are exactly loop, path-recip, 2-recip, and 3-recip-triangles. (These are given in light blue, yellow, brown, and pink, respectively.) Now consider transitive relations that are *not* reciprocated. For example, A connects to B who connects to C , but A does not connect to C . When a triangle contains a cycle, a reciprocated transitive relationship creates a reciprocal edges.

Since loop-triangles have no reciprocal edges, there are three transitive relations that are not reciprocated. Analogously, for path-recip-triangle, there are two such unreciprocated relations. And for 2-recip and 3-recip triangles, these numbers are one and zero.

So we ask, when a triangle contains a cycle, does it contain unreciprocated transitive relations? One would think that a cycle indicates a strong tie between three vertices,

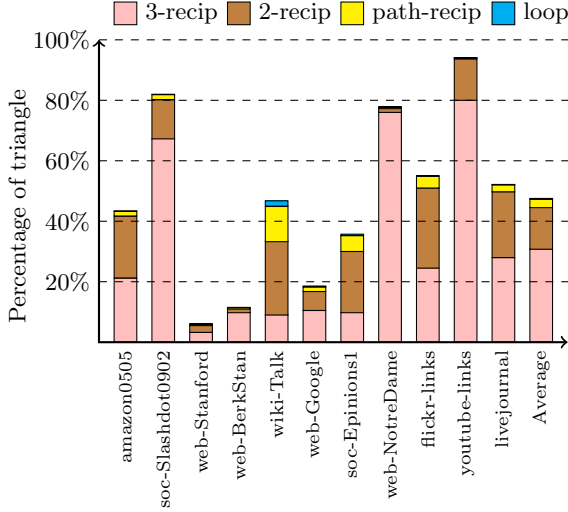


Figure 9: Power of transitivity: For each graph in our collection, we plot the percentages of (different) triangle types containing a cycle. Each bar corresponds to a single graph, and the stacked bar charts gives the percentages of the 4 different triangle types. Note the dominance of pink and brown (3-recip and 2-recip-triangles).

$\frac{(1-r)^2}{4}$	$\frac{(1-r)^2}{2}$	$\frac{(1-r)^2}{4}$	$r(1-r)$	$r(1-r)$	r^2

Table 4: The probability of an undirected wedge become a particular directed type wedge

and so reciprocation is expected. This is exactly what we see in Figure 9, quite strongly over practically all graphs. Almost all triangles with a cycle are either 2-recip or 3-recip-triangles. We almost never see any loop-triangles, shown by the lack of light blue in Figure 9. Again, this is more evidence that reciprocal edges play an important role in graph structure. The results demonstrate that the power of transitivity of real world networks. One can observe that social relationships carried forward two steps (as a transitive relation) almost always lead to reciprocation.

4. NULL MODELS FOR (ψ, τ) -CLOSURE

In the previous section, we made several observations about the (ψ, τ) -closure rates in real graphs. How significant are these results? Can they be explained merely by the reciprocity of a graph? We propose a *null hypothesis*, based on assigning the type of each edge only based on the reciprocity of the graph. We start by making the graph undirected and insert direction and reciprocity randomly as follows. If (u, v) is an undirected edge, we make it reciprocal with probability r ; we direct it from u to v with probability $(1-r)/2$, we direct from v to u with probability $(1-r)/2$. Based on this model, the probabilities of an undirected wedge and/or triangle being of a certain type can be calculated through simple calculations. This information is presented in Table 4 and Table 5.

$\frac{3(1-r)^3}{4}$	$\frac{(1-r)^3}{4}$	$\frac{3r(1-r)^2}{4}$	$\frac{3r(1-r)^2}{2}$
$\frac{3r(1-r)^2}{4}$	$3r^2(1-r)$	r^3	

Table 5: Fractions of triangle types based on the null model.

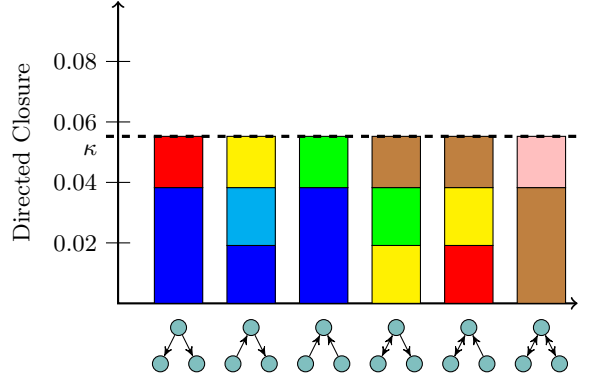


Figure 10: Closure chart of web-Google for random directions: We consider the undirected version of web-Google graph and added one-way and reciprocal edges according the random null model. Observe that the total closure for each wedge is identical, and how different this is from Figure 2.

Table 5 reveals that observations of the previous section cannot be explained by randomness or reciprocity. For instance, if we compare the expected fractions of the last two triangles 2-recip and 3-recip, we see that 2-recip should be more frequent when the reciprocity, $r < 0.75$. Even though this condition holds in most of the graphs in our data set, we observe the contrary behavior in real data sets, and 3-recip generally is more frequent than 2-recip. Another observation is about loop triangles. According to our null model, trans and loop triangles have the same dependence of reciprocity, and trans triangles are expected to be only 3 times more frequent than loop triangles. However, transtriangles are much more frequent in practice. In other words, the null model can explain the sparsity of looptriangles, but not their near absence.

Figure 10 illustrates how the directed closure chart would look when direction is random. We take the undirected version of web-Google graph and add one-way and reciprocal edges according the random null model. If we compare this figure to Figure 2, we see a totally different distribution, pointing to the significance of our results. Here, we are only presenting the results for web-Google due to space limitations, but we observed the same trend in all other graphs.

Finally, in Figure 11 we look at two triangle types, out-recip and path-recip, whose dependences on reciprocity are the same, but one is overrepresented, while the other is under represented, compared to the expectation of the null hypothesis. Type out-recip-triangles are overrepresented in

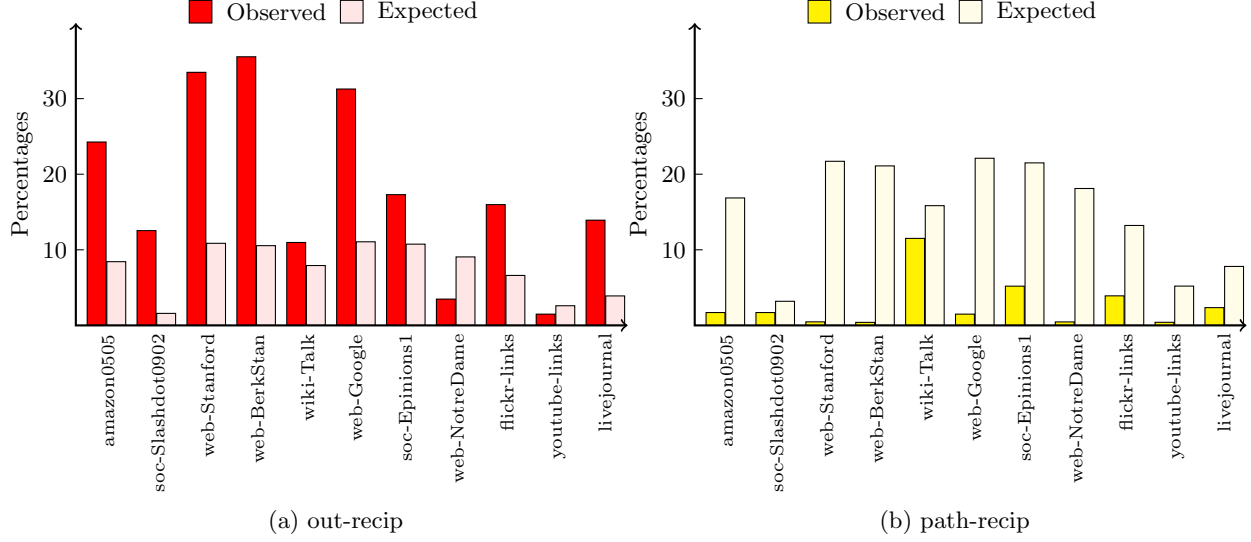


Figure 11: Deviations from the null model: For various graphs, we plot the fraction of triangles of a given type together with what is predicted by the random null model. This is done for the out-recipe and path-recipe-triangles. Observe the large differences, showing that directed triangle distributions are far from random.

all graphs except **web-NotreDame** and **youtube-links**, while path-recipe-triangles are underrepresented in all graphs.

All these results show that the direction in triangles reveals a special structure, which cannot be explained by randomness or reciprocity.

5. COUNTING DIRECTED TRIANGLES

The results in the previous section showed the importance of computing directed closure charts. In this section, we turn our attention how to perform this task efficiently and describe approximation algorithms to estimate the various clustering coefficients (and also the numbers of triangles). We extend the method of *wedge sampling* for directed graphs.

We begin with some basic notation. We define the following seven subsets of W_ψ . Let

$$W_\psi(\tau) = \{ w \in W_\psi \mid w \text{ is } \tau\text{-closed} \}.$$

Note that $\kappa_{\psi,\tau} = |W_\psi(\tau)|/|W_\psi|$. This fraction can now be estimated through the following algorithmic template.

1. Select k uniform random ψ -wedges (with replacement).
2. Determine k' , the number of τ -closed wedges among this sample.
3. Output estimate $\hat{\kappa}_{\psi,\tau} = k'/k$ for $\kappa_{\psi,\tau}$.

The main theorem shows that this provides a good estimate for $\kappa_{\psi,\tau}$. Similar versions of this theorem have appeared in our earlier work [19, 20], but we provide a proof for completeness. We first state Hoeffding's inequality.

THEOREM 5.1 (HOEFFDING [12]). *Let X_1, X_2, \dots, X_k be independent random variables with $0 \leq X_i \leq 1$ for all $i = 1, \dots, k$. Define $\bar{X} = \frac{1}{k} \sum_{i=1}^k X_i$. Then for any $\varepsilon > 0$, we have*

$$\Pr[|\bar{X} - \mathbb{E}[\bar{X}]| \geq t] \leq 2 \exp(-2t^2/k).$$

THEOREM 5.2. *Let $\varepsilon, \delta > 0$ and set $k = \lceil 0.5 \varepsilon^{-2} \ln(2/\delta) \rceil$.*

$$\Pr[|\hat{\kappa}_{\psi,\tau} - \kappa_{\psi,\tau}| \geq \varepsilon] \leq \delta$$

PROOF. Define indicator random variable X_i for the i th ψ -wedge being τ -closed (so $X_i = 1$ if the wedge is τ -closed and 0 otherwise). Note that $\mathbb{E}[X_i] = \kappa_{\psi,\tau}$, so $\mathbb{E}[\sum_{i=1}^k X_i] = k\kappa_{\psi,\tau}$. Since $\hat{\kappa}_{\psi,\tau} = \sum_{i=1}^k X_i/k$, the event $|\hat{\kappa}_{\psi,\tau} - \kappa_{\psi,\tau}| \geq \varepsilon$ is the same as $|\sum_{i=1}^k X_i - \mathbb{E}[\sum_{i=1}^k X_i]| \geq \varepsilon k$. By [Theorem 5.1](#), the probability of this event, by choice of k , is at most $2 \exp(-2\varepsilon^2 k^2/k) < \delta$. \square

A direct corollary of this theorem gives bounds for triangle estimates. This is obtained by multiplying event inequality in [Theorem 5.2](#) by $|W_\psi|/\chi(\psi, \tau)$ and observing that $|T_\tau| = \kappa_{\psi,\tau}|W_\psi|/\chi(\psi, \tau)$.

COROLLARY 5.3. *Fix types ψ and τ such that $\chi(\psi, \tau) \neq 0$. Denote $\hat{T} = \hat{\kappa}_{\psi,\tau}|W_\psi|/\chi(\psi, \tau)$.*

$$\Pr[|\hat{T} - |T_\tau|| \geq \varepsilon|W_\psi|/\chi(\psi, \tau)] \leq \delta$$

There are a few subtleties here worth mentioning. For the same number of samples, we can use *different* wedge types to count the same triangle set $|T_\tau|$. For a candidate type ψ , the error is proportional to $|W_\psi|/\chi(\psi, \tau)$. Hence, using wedge types that are less frequent give stronger approximations for the same triangle type. Another consequence of this observation is that only 4 wedge types (e.g., in, recip-out, recip-in, and recip-tot) are sufficient to compute the numbers of all triangles types and thus the 15 closure rates.

5.1 Uniform sampling of wedge types

To give a full algorithm, we need to give a procedure that samples uniform random wedges of any desired type. We can split wedge types into two groups: homogenous and heterogenous. Homogenous wedges have only one kind of edge, such as in, out, and recip-tot-wedges. Heterogenous wedges have different kinds of wedges, such as mid, recip-in,

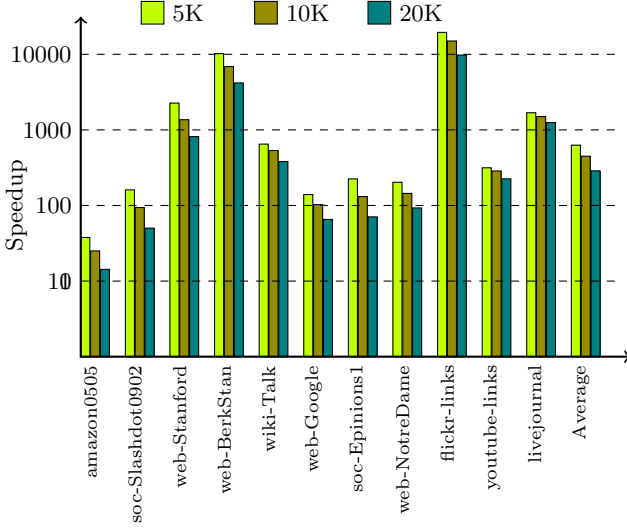


Figure 12: Speed-up over enumeration: We count the total number of directed triangles through wedge sampling and compare it to enumeration methods. We runs our algorithm for 5K, 10K, and 20K samples.

and recip-out-wedges. The sampling procedures are analogous for types within a group. Hence, we will only describe how to sample uniform random in-wedges and random mid-wedges.

First, we deal with in-wedges. Set $p_v = \binom{d_v^+}{2} / |W_{ii}|$, where $v \in V$. Note that $\sum_{v \in V} p_v = 1$, so this forms a probability distribution over V .

- Sample a random v according to the distribution given by $\{p_v\}$.
- Sample a uniform random pair u, w of in-neighbors of v .
- Output the wedge $\{(u, v), (v, w)\}$

This generates a uniform random in-wedge. The number of in-wedges incident to v is exactly $p_v |W_{ii}|$, and the second step generates a uniform random in-wedge centered at v .

Now for out-wedges. Set $p_v = d_v^+ d_v^- / |W_{ii}|$, where $v \in V$. Again, $\sum_{v \in V} p_v = 1$.

- Sample a random v according to the distribution given by $\{p_v\}$.
- Sample u , a uniform random in-neighbor of v , and w , a uniform random out-neighbor.
- Output the wedge $\{(u, v), (v, w)\}$

We can show that is a uniform random out-wedge, using an argument almost identical to that used above.

With these procedures, we can implement the wedge sampling algorithms for all wedge/triangle types.

5.2 Experimental Results

We implemented our algorithms in C and ran our experiments on a computer equipped with a 2.3GHz Intel core i7 processor with 4 cores and 256KB L2 cache (per core), 8MB L3 cache, and 8GB memory. We performed our experiments on 11 graphs, whose properties are presented in Table 3.

In Figure 12, we compare the runtime of wedge sampling to the best enumeration algorithm. Our enumeration algorithm is based on the principles of [4, 19, 6, 24], such that

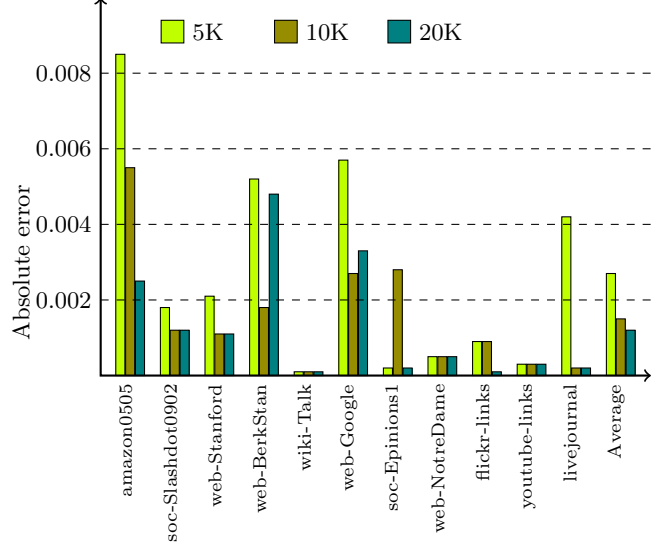


Figure 13: Improved accuracy with more wedge samples: We focus on estimating the fraction that a out-wedge closed to a out-recipe-triangle ($\kappa_{iii,c}$). We consider 5K, 10K, 20K wedge samples. The errors are all in third decimal point.

each edge is assigned to its vertex with a smaller total degree, d_v , (using the vertex numbering as a tie-breaker), and then vertices only check closure for wedges formed by edges assigned to them. Once a triangle is identified, it is classified according to its edges. As seen in Figure 12, wedge sampling works orders of magnitude faster than the enumeration algorithm. The timing results show tremendous savings; for instance, wedge sampling only takes 0.064 seconds on **web-BerkStan** while full enumeration takes 271 seconds.

Figure 13 shows the accuracy of the wedge sampling algorithm, by displaying the sampling error in computing how often a out-wedge closes to a out-recipe-triangle. At 99.9% confidence ($\delta = 0.001$), the upper bound on the error we expect for 5K, 10K, and 20K samples is .028, .020, and .013, respectively. In all our experiments, the observed error is always much smaller than what is indicated by Theorem 5.1. For instance, the maximum error for 5K samples is .0085, much less than that 0.028 given by the upper bound.

Due to space limitations, we cannot present results accuracies for other closer rates. However, we can report that the proposed wedge sampling algorithm produced consistently more accurate results than what is indicated by Theorem 5.1 for all graphs and all closure rates.

6. CONCLUSIONS

We initiate the study of directed triangles in massive networks, by defining the set of directed closure measures. These quantities reveal a surprising amount of information about digraphs. He observe heterogeneity in closure rates of different wedges, strong effect of reciprocity in closure rates, the power of transitivity in the structure of triangles. Our results also show that these observations cannot be explained merely by randomness or reciprocity. We hope that this paper leads the way in deeper studies into digraphs, and also convinces the data mining and social networks community that direction cannot be ignored. The fast estimation re-

sults show that the measures can be computed in a scalable manner.

7. REFERENCES

- [1] Laboratory for web algorithmics. Available at <http://law.di.unimi.it/webdata/ljournal-2008/>.
- [2] S. E. Ahnert and T. M. A. Fink. Clustering signatures classify directed networks. *Phys. Rev. E*, 78:036112, Sep 2008.
- [3] D. Cartwright and F. Harary. Structural balance: a generalization of heider’s theory. *Psychological Review*, 63(5):277–293, 1956.
- [4] N. Chiba and T. Nishizeki. Arboricity and subgraph listing algorithms. *SIAM J. Comput.*, 14:210–223, February 1985.
- [5] F. Chierichetti, R. Kumar, S. Lattanzi, M. Mitzenmacher, A. Panconesi, and P. Raghavan. On compressing social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 219–288, 2009.
- [6] J. Cohen. Graph twiddling in a MapReduce world. *Computing in Science & Engineering*, 11:29–41, 2009.
- [7] N. Durak, T. G. Kolda, A. Pinar, and C. Seshadhri. A scalable directed graph model with reciprocal edges. In *Proc. IEEE 2nd Int. Workshop on Network Science*, 2013. also available as arxiv:1210.5288.
- [8] G. Fagiolo. Clustering in complex directed networks. *Phys. Rev. E*, 76:026107, Aug 2007.
- [9] K. Faust. Comparing social networks: Size, density, and local structure. *Metodoloski zvezki*, 3(2):185–216, 2006.
- [10] D. Garlaschelli and M. I. Loffredo. Patterns of link reciprocity in directed networks. *Phys. Rev. Lett.*, 93:268701, Dec. 2004.
- [11] F. Heider. Attitudes and cognitive organization. *Journal of Psychology*, 21:107–112, 1946.
- [12] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- [13] P. W. Holland and S. Leinhardt. A method for detecting structure in sociometric data. *American Journal of Sociology*, 76:492–513, 1970.
- [14] H. Kwak, C. Lee, H. Park, and S. Moon. What is twitter, a social network or a news media? In *WWW ’10*, pages 591–600. ACM, 2010.
- [15] J. Leskovec, D. Huttenlocher, and J. Kleinberg. Signed networks in social media. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1361–1370, 2010.
- [16] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network motifs: Simple building blocks of complex networks. *Science*, 298(5594):824–827.
- [17] A. Mislove, H. S. Koppula, K. P. Gummadi, P. Druschel, and B. Bhattacharjee. Growth of the flickr social network. In *WOSN’08*, pages 25–30. ACM, 2008.
- [18] M. E. J. Newman, S. Forrest, and J. Balthrop. Email networks and the spread of computer viruses. *Phys. Rev. E*, 66(3):035101, Sept. 2002.
- [19] T. Schank and D. Wagner. Finding, counting and listing all triangles in large graphs, an experimental study. In *Experimental and Efficient Algorithms*, pages 606–609. Springer Berlin / Heidelberg, 2005.
- [20] C. Seshadhri, A. Pinar, and T. G. Kolda. Triadic measures on graphs: The power of wedge sampling. In *Proceedings of the SIAM Conference on Data Mining (SDM)*, 2013.
- [21] J. Skvoretz. Biased net theory: Approximations, simulations and observations. *Social N*, 12(3):217–238, 1990.
- [22] J. Skvoretz, T. J. Fararo, and F. Agneessens. Advances in biased net theory: definitions, derivations, and estimations. *Social Networks*, 26:113–119, 2004.
- [23] S. Son, A. R. Kang, H.-c. Kim, T. Kwon, J. Park, and H. K. Kim. Analysis of context dependence in social interaction networks of a massively multiplayer online role-playing game. *PLoS ONE*, 7(4):e33918, 04 2012.
- [24] S. Suri and S. Vassilvitskii. Counting triangles and the curse of the last reducer. In *WWW’11*, pages 607–614, 2011.
- [25] S. Wasserman and K. Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, 1994.
- [26] D. Watts and S. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393:440–442, 1998.
- [27] Stanford Network Analysis Project (SNAP). Available at <http://snap.stanford.edu/>.